

# RASA: RANK-SHARING LOW-RANK ADAPTATION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Low-rank adaptation (LoRA) has been prominently employed for parameter-efficient fine-tuning of large language models (LLMs). However, the limited expressive capacity of LoRA, stemming from the low-rank constraint, has been recognized as a bottleneck, particularly in rigorous tasks like code generation and mathematical reasoning. To address this limitation, we introduce Rank-Sharing Low-Rank Adaptation (RaSA), an innovative extension that enhances the expressive capacity of LoRA by leveraging partial rank sharing across layers. By forming a shared rank pool and applying layer-specific weighting, RaSA effectively increases the number of ranks without augmenting parameter overhead. Our theoretically grounded and empirically validated approach demonstrates that RaSA not only maintains the core advantages of LoRA but also significantly boosts performance in challenging code and math tasks. Code, data and scripts are available at: <https://anonymous.4open.science/r/RaSA-ICLR-0E25>.

## 1 INTRODUCTION

Low-rank adaptation (LoRA, Hu et al. (2022)) has become a de facto parameter-efficient fine-tuning (PEFT) method for adapting large language models (LLMs) to specific downstream tasks. Its core idea is to constrain the parameter updates to be low-rank, which significantly reduces the number of trainable parameters and allows them to be merged back into the original model, thereby avoiding additional inference latency. Despite its advantages, recent studies have shown that LoRA still lags behind full fine-tuning (FFT), particularly in scenarios involving large training datasets and complex tasks such as mathematical reasoning and code generation (Jiang et al., 2024; Biderman et al., 2024). A plausible explanation for this performance gap is that the low-rank constraint limits the expressive capacity of LoRA. For instance, Biderman et al. (2024) empirically found that the effective rank required for FFT is 10-100× higher than typical LoRA configuration, and Zeng & Lee (2024) theoretically demonstrated that a Transformer network (Vaswani et al., 2017) requires a rank at least half the size of the model dimension to approximate another model of similar size.

Although the limited number of trainable parameters results in limited expressive capacity, recent studies still indicate redundancy in LoRA’s parameters. For example, Kopiczko et al. (2024); Song et al. (2024); Renduchintala et al. (2024); Li et al. (2024) further reduced the number of LoRA’s parameters by sharing them across layers and modules with only slight performance loss. Brühl-Gabrielsson et al. (2024) compressed 1,000 LoRAs trained from different tasks by sharing their parameter spaces. This contradiction suggests that LoRA’s parameters are still not being fully utilized.

Combining the above two observations, we propose **Rank-Sharing Low-Rank Adaptation (RaSA)**, an approach that boosts the expressive capacity of LoRA by enabling partial rank sharing across layers. Specifically, given an LLM with  $L$  layers, RaSA extracts  $k$  ranks from each layer’s LoRA update to form a rank pool of  $L \times k$  ranks, which is shared across all layers with layer-specific weighting. RaSA retains the core advantages of LoRA – keeping the same parameter overhead and allowing for easy merging back into the model. Moreover, since modern LLMs typically have deep architectures (i.e., large  $L$ ), RaSA greatly increase the effective rank of the parameter update by  $(L - 1) \times k$ .

However, a higher rank does not necessarily lead to better expressive capacity. To rigorously assess the benefits of RaSA, we analyze its capacity to reconstruct high-rank matrices compared to LoRA. Theoretically, we prove that RaSA’s minimum reconstruction error is bounded by that of LoRA. Empirically, we show that when  $k$  is relatively small, RaSA can be easily optimized to achieve a significantly lower reconstruction error than LoRA. Finally, we conducted experiments

on mathematical reasoning and code generation, demonstrating that the lower reconstruction error translates to improved downstream task performance.

Our contributions are summarized as follows:

- We propose RaSA, a novel extension of LoRA by allowing partial rank sharing across layers, which significantly improves the method’s efficiency and expressiveness (§ 2).
- We provide a comprehensive analysis – both theoretical and empirical – showcasing RaSA’s superior capacity for matrix reconstruction (§ 3) and its resultant improved performance on rigorous downstream tasks (e.g. code and math) (§ 4).

## 2 METHOD

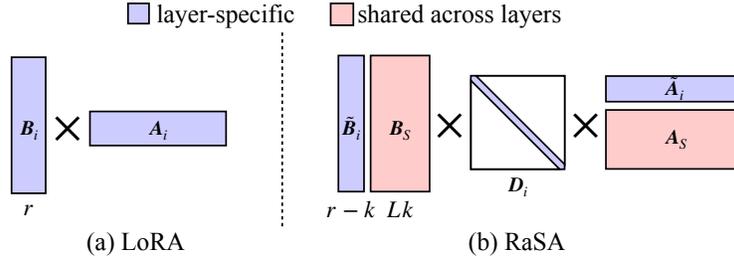


Figure 1: Decomposition of the update matrix  $\Delta W_i$  in LoRA and RaSA, where  $i$  is the layer index.

### 2.1 FORMULATION

Given a pre-trained weight matrix  $W \in \mathbb{R}^{b \times a}$ , LoRA constrains its update to a low-rank form by decomposing the update matrix  $\Delta W \in \mathbb{R}^{b \times a}$  into a product of two rank- $r$  matrices:

$$W + \Delta W = W + \frac{\alpha}{r} BA \quad (B \in \mathbb{R}^{b \times r}, A \in \mathbb{R}^{r \times a}), \quad (1)$$

where rank  $r \ll \min(b, a)$  serves as a bottleneck dimension, reducing the number of trainable parameters, and  $\alpha$  is a scaling factor. In an LLM with  $L$  layers, LoRA assigns distinct trainable matrices to each layer- $i$ :  $\{B_i, A_i\}_{i \in [L]}$  (Figure 1(a)). RaSA, on the other hand, mitigates the low-rank bottleneck of LoRA through rank sharing. Specifically, RaSA takes out  $k$  ranks in each layer and shares them across all layers. This process can be conceptualized as follows:

1. Split the matrices  $B_i$  and  $A_i$  into layer-specific parts ( $\tilde{B}_i, \tilde{A}_i$ ) and layer-shared parts ( $\hat{B}_i, \hat{A}_i$ ):

$$B_i = \begin{bmatrix} \underbrace{\tilde{B}_i}_{\mathbb{R}^{b \times (r-k)}} & \underbrace{\hat{B}_i}_{\mathbb{R}^{b \times k}} \end{bmatrix}, \quad A_i = \begin{bmatrix} \underbrace{\tilde{A}_i^T}_{\mathbb{R}^{a \times (r-k)}} & \underbrace{\hat{A}_i^T}_{\mathbb{R}^{a \times k}} \end{bmatrix}^T. \quad (2)$$

2. Concatenate the layer-shared parts across all layers to form shared rank pools ( $B_S$  and  $A_S$ ):

$$B_S = [\hat{B}_1 \quad \hat{B}_2 \quad \dots \quad \hat{B}_L] \in \mathbb{R}^{b \times (L \times k)}, \quad A_S = [\hat{A}_1^T \quad \hat{A}_2^T \quad \dots \quad \hat{A}_L^T]^T \in \mathbb{R}^{(L \times k) \times a}. \quad (3)$$

Therefore, the update for layer- $i$  is given by:

$$\begin{aligned} W_i + \Delta W_i &= W_i + \frac{\alpha}{r} (\tilde{B}_i \tilde{A}_i + B_S A_S) \\ &= W_i + [\tilde{B}_i \quad B_S] \text{diag}\left(\frac{\alpha}{r}\right) \begin{bmatrix} \tilde{A}_i \\ A_S \end{bmatrix}. \end{aligned} \quad (4)$$

To enable layer-specific weighting, we replace the constant diagonal matrix with a trainable diagonal matrix  $D_i = \text{diag}(d_1, d_2, \dots, d_j, \dots, d_{r-k+Lk})$ , yielding the final RaSA update (Figure 1(b)):

$$W_i + \Delta W_i = W_i + \underbrace{[\tilde{B}_i \quad B_S]}_{\mathbb{R}^{b \times (r-k+Lk)}} D_i \underbrace{\begin{bmatrix} \tilde{A}_i \\ A_S \end{bmatrix}}_{\mathbb{R}^{(r-k+Lk) \times a}}. \quad (5)$$

## 2.2 ANALYSIS & IMPLEMENTATION DETAILS

**Rank of  $\Delta W$**  Comparing Equations (1) and (5), RaSA increases the rank of  $\Delta W$  from  $r$  to  $r - k + Lk$ . Since modern LLMs are deep, RaSA significantly boosts the model’s expressive capacity by enabling a higher effective rank, on which we have a detailed discussion in § 3. [Each layer in RaSA maintains the same rank for  \$\Delta W\$ , which sets it apart from methods that dynamically assign ranks across layers, such as AdaLoRA \(Zhang et al., 2023\) and PriLoRA \(Benedek & Wolf, 2024\).](#)

**Additional Parameters** RaSA introduces the diagonal matrix  $D_i$  as additional parameters. Since  $D_i$  is diagonal and operates only at the bottleneck dimension, the added parameters are negligible. In practice,  $D_i$  contributes to less than 0.001% of the total model parameters.

**Initialization** Following LoRA, we use Kaiming initialization (He et al., 2015) for  $\tilde{A}_i$  and  $A_S$ , and initialize  $\tilde{B}_i$  and  $B_S$  to zero. For  $D_i$ , we differentiate between the layer-specific and layer-shared parts by scaling  $\alpha$  proportionally by their respective ranks:

$$d_j = \begin{cases} \frac{1}{2} \times \frac{\alpha}{r-k} & \text{if } j \leq r - k, \\ \frac{1}{2} \times \frac{\alpha}{Lk} & \text{if } j > r - k. \end{cases} \quad (6)$$

**Same Dimension Assumption** RaSA assumes that all layers share the same dimensionality. This holds for the vast majority of models (e.g. Llama (Dubey et al., 2024), Mistral (Jiang et al., 2023)).

## 3 RECONSTRUCTION ERROR ANALYSIS

While RaSA increases the effective rank of  $\Delta W$ , a higher rank does not necessarily guarantee improved expressive capacity. For instance, a full-rank identity matrix can only perform the identity transformation. To assess the expressive capacity of LoRA and RaSA, we compare their abilities to reconstruct a set of high-rank matrices  $\{M_i\}_{i \in [L]}$ , where  $\text{rank}(M_i) = R > r$ . Under the Frobenius norm, the **minimum reconstruction error (MRE) of LoRA** is defined as:

$$e_{\text{loara}} = \min_{B_i, A_i} \sum_{i=1}^L \|M_i - B_i A_i\|_F^2. \quad (7)$$

According to the Eckart–Young–Mirsky theorem (Eckart & Young, 1936), we can perform singular value decomposition (SVD) on  $M_i$ :

$$\text{SVD}(M_i) = \sum_{j=1}^R \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} \quad (\sigma_1^{(i)} \geq \sigma_2^{(i)} \geq \dots \geq \sigma_R^{(i)}). \quad (8)$$

LoRA’s optimal approximation is given by the first  $r$  components of  $\text{SVD}(M_i)$ , and  $e_{\text{loara}}$  becomes the sum of squares of the discarded singular values (those beyond the  $r$ -th one):

$$e_{\text{loara}} = \sum_{i=1}^L \|M_i - \sum_{j=1}^r \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T}\|_F^2 = \sum_{i=1}^L \sum_{j=r+1}^R \sigma_j^{(i)2}. \quad (9)$$

Similarly, when each layer shares  $k$  ranks out, we can define the **MRE of RaSA** as:

$$e_{\text{rasa}(k)} = \min_{\tilde{B}_i, \tilde{A}_i, B_S, A_S, D_i} \sum_{i=1}^L \|M_i - [\tilde{B}_i \ B_S] D_i \begin{bmatrix} \tilde{A}_i \\ A_S \end{bmatrix}\|_F^2. \quad (10)$$

For simplicity, in this section we consider that  $D_i$  operates only on the shared matrices  $B_S$  and  $A_S$ , which does not affect the value of  $e_{\text{rasa}(k)}$ :

$$e_{\text{rasa}(k)} = \min_{\tilde{B}_i, \tilde{A}_i, B_S, A_S, D_i} \sum_{i=1}^L \|M_i - (\tilde{B}_i \tilde{A}_i + B_S D_i A_S)\|_F^2. \quad (11)$$

### 3.1 THEORETICAL ANALYSIS

**Theorem 3.1.**  $e_{\text{rasa}(k)} \leq e_{\text{loara}}$

*Proof.* To prove this, we construct a feasible solution for RaSA that achieves the same reconstruction error as LoRA’s minimum error. This is done by distributing the ranks shared across layers in RaSA such that they cover the same rank range as the optimal LoRA solution.

For each layer- $i$ , we take the last  $k$  components (corresponding to the indices  $r - k + 1$  through  $r$ ) from the LoRA’s optimal approximation (Equation (9)), forming the following matrices:

$$\begin{aligned} \mathbf{U}^{(i)} &= [\mathbf{u}_{r-k+1}^{(i)} \quad \mathbf{u}_{r-k+2}^{(i)} \quad \cdots \quad \mathbf{u}_r^{(i)}], \\ \mathbf{V}^{(i)} &= [\mathbf{v}_{r-k+1}^{(i)} \quad \mathbf{v}_{r-k+2}^{(i)} \quad \cdots \quad \mathbf{v}_r^{(i)}], \\ \Sigma^{(i)} &= [\sigma_{r-k+1}^{(i)} \quad \sigma_{r-k+2}^{(i)} \quad \cdots \quad \sigma_r^{(i)}]. \end{aligned} \quad (12)$$

The shared matrices  $\mathbf{B}_S$  and  $\mathbf{A}_S$  are constructed by stacking  $\mathbf{U}^{(i)}$  and  $\mathbf{V}^{(i)}$  from each layer:

$$\mathbf{B}_S = [\mathbf{U}^{(1)} \quad \cdots \quad \mathbf{U}^{(i)} \quad \cdots \quad \mathbf{U}^{(L)}], \quad \mathbf{A}_S = [\mathbf{V}^{(1)} \quad \cdots \quad \mathbf{V}^{(i)} \quad \cdots \quad \mathbf{V}^{(L)}]^T. \quad (13)$$

Similarly, we define the diagonal matrix  $\mathbf{D}_i$  for each layer- $i$  by placing the corresponding singular values  $\Sigma^{(i)}$  in their appropriate positions:

$$\mathbf{D}_i = \text{diag}([\mathbf{0} \quad \cdots \quad \Sigma^{(i)} \quad \cdots \quad \mathbf{0}]). \quad (14)$$

Finally, the matrices  $\tilde{\mathbf{B}}_i$  and  $\tilde{\mathbf{A}}_i$  are formed from the first  $r - k$  components of  $\text{SVD}(\mathbf{M}_i)$ :

$$\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i = \sum_{j=1}^{r-k} \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T}. \quad (15)$$

Substituting Equations (13) to (15) into Equation (11), we derive the following:

$$\begin{aligned} e_{\text{rasa}(k)} &\leq \sum_i^L \|\mathbf{M}_i - (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S)\|_F^2 \\ &= \sum_{i=1}^L \left\| \sum_{j=1}^R \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} - \left( \sum_{j=1}^{r-k} \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} + \sum_{j=r-k+1}^r \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} \right) \right\|_F^2 \\ &= \sum_{i=1}^L \left\| \sum_{j=1}^R \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} - \sum_{j=1}^r \sigma_j^{(i)} \mathbf{u}_j^{(i)} \mathbf{v}_j^{(i)T} \right\|_F^2 \\ &= \sum_{i=1}^L \sum_{j=r+1}^R \sigma_j^{(i)2} \\ &= e_{\text{loa}}. \end{aligned} \quad (16)$$

Thus, we conclude that  $e_{\text{rasa}(k)} \leq e_{\text{loa}}$ , proving that RaSA can achieve equal or lower minimum reconstruction error compared to LoRA.  $\square$

### 3.2 EMPIRICAL ANALYSIS

While the previous theoretical analysis guarantees that RaSA can at least match the MRE of LoRA, it does not quantify how much RaSA improves upon LoRA. To provide a more intuitive understanding of how RaSA achieves lower reconstruction error, we turn to an optimization-based analysis using coordinate descent.

**Empirical Validation** Specifically, we instantiate the set of high-rank matrices  $\{\mathbf{M}_i\}_{i \in [L]}$  with the actual weight updates from model fine-tuning:  $\{\Delta \mathbf{W}_i\}_{i \in [L]}$ , and iteratively minimize the reconstruction error in Equation (11) by adjusting the parameters of RaSA ( $r = 8, k = 1$ ), namely the  $\tilde{\mathbf{B}}_i$ ,  $\tilde{\mathbf{A}}_i$ ,  $\mathbf{B}_S$ ,  $\mathbf{A}_S$ , and  $\mathbf{D}_i$  (details can be found in appendix A). We apply this procedure to various kinds of linear modules within Llama-3.1-8B until convergence, and compute  $e_{\text{loa}}$  using Equation (9) as baseline values.

Figure 2 shows that RaSA requires  $\sim 10$  iterations to achieve a significantly lower reconstruction error than LoRA’s minimum. This pattern is consistent across all linear modules in the model, demonstrating the enhanced expressive capacity of RaSA.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

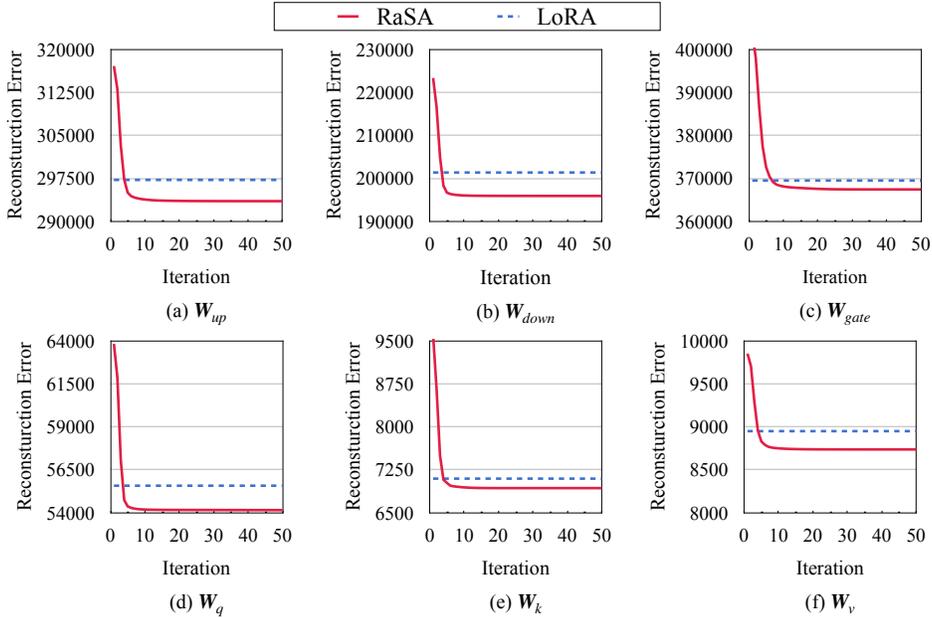


Figure 2: Reconstruction error curves of RaSA ( $r = 8, k = 1$ ) during coordinate descent. We also plot the minimum reconstruction error of LoRA (Equation (9)) for comparison.

**Selection of  $k$**  RaSA introduces only one additional hyper-parameter,  $k$ , which controls how many ranks are taken from each layer to be shared across all layers. When  $k = 0$ , RaSA reduces to LoRA, where no ranks are shared. On the other hand, when  $k = r$ , RaSA shares all ranks across layers, eliminating layer-specific low-rank updates and making the adaptation fully shared. While this maximizes the effective rank of update, it may diminish layer diversity and the ability to capture layer-specific nuances. We traversed  $k$  from 0 to 8, and presents the converged reconstruction error from the previous coordinate descent experiment in Figure 3. The results indicate that a small value of  $k$ , around  $r/8$ , achieves the minimum error. Further increasing  $k$  can lead to a rise in reconstruction error, even exceeding that of LoRA. This finding also indicates that some current methods that share all ranks across all layers, such as VeRA (Kopiczko et al., 2024) and Tied-LoRA (Renduchintala et al., 2024), might be sub-optimal and challenging to be optimized. *It is worth noting that  $k = r/8$  is still an empirical result. We leave the theoretically optimal  $k$  to future work, which might also be related the model dimension and number of layers.*

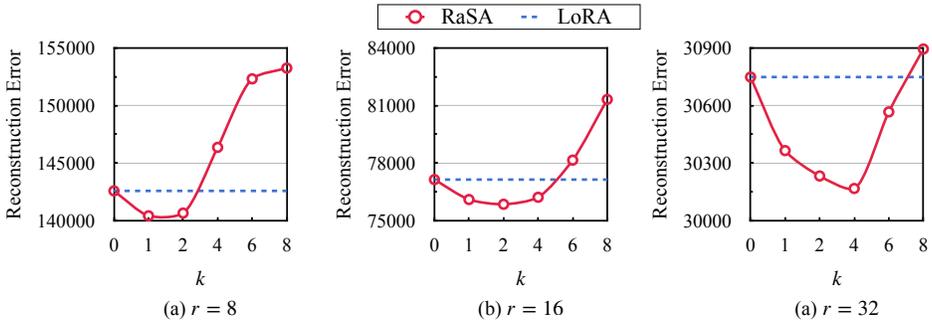


Figure 3: Reconstruction error comparison between RaSA and LoRA as a function of the shared rank parameter  $k$ . We also plot the minimum reconstruction error of LoRA (Equation (9)) for comparison. The results are average across all linear modules in the model.

## 4 EXPERIMENT

### 4.1 SETUP

**Tasks** Our experiments generally align with those reported by Biderman et al. (2024). We applied all the methods to instruction fine-tuning and evaluated their performance on challenging tasks: code generation and mathematical reasoning. While Biderman et al. (2024) use Humaneval (Chen et al., 2021) and GSM8K (Cobbe et al., 2021) as test sets, these two benchmarks have become saturated with the rapid growth of LLMs. To provide a more rigorous evaluation, we adopted two more challenging benchmarks as test sets. Since these two benchmarks lack validation sets, in addition to reporting the results from the last checkpoint, we also report the best results as a reference for the upper bound of each method. Prompt templates for evaluation are provided in appendix B.

- *Code Generation*: We used Magicoder-Evol-Instruct-110k (Wei et al., 2024) as the training data, a collection of programming question-answer (QA) pairs, which is a reproduced and decontaminated version of WizardCoder (Luo et al., 2024). We used Humaneval+ (Liu et al., 2023) as the test set, an extension of the Humaneval benchmark that scales the number of test cases by  $80\times$ . We used the Bigcode Evaluation Harness (Ben Allal et al., 2022) as the evaluation tool, sampling 50 solutions per problem with a temperature of 0.2, and report both Pass@1 and Pass@10.
- *Mathematical Reasoning*: We used MetaMathQA (Yu et al., 2024) as the training data, which comprises 395K QA pairs derived from the training sets of GSM8K (Cobbe et al., 2021) and MATH (Hendrycks et al., 2021), rewritten by GPT-3.5. We used MATH (Hendrycks et al., 2021) as the test set, which consists of 5K competition-level mathematics problems covering 7 subjects and 5 difficulty levels. We followed the evaluation protocol from LLMs Evaluation Harness (Gao et al., 2024), using sympy to verify correctness and employing greedy search for generation.

**Baselines** We compare RaSA to the several representative PEFT methods:

- LoRA (Hu et al., 2022) that learns only a low-rank perturbation to the pretrained weight matrix.
- MoRA (Jiang et al., 2024) that uses block diagonal matrices instead of low-rank matrices.
- VeRA (Kopiczko et al., 2024) that fully shares the low-rank matrices across all layers with layer-specific weighting, and freeze the low-rank matrices during training to achieve extreme parameter efficient fine-tuning. Therefore, VeRA can set a higher rank- $r$  than LoRA.

**LLMs & Training Details** We conducted experiments on two open-sourced LLMs: Llama-3.1-8B (Dubey et al., 2024) and Mistral-0.3-7B (Jiang et al., 2023). Following common practice (Kopiczko et al., 2024; Jiang et al., 2024), we used pre-trained models rather than instruction-tuned ones. We applied PEFTs on all linear modules from attention ( $\mathbf{W}_q, \mathbf{W}_k, \mathbf{W}_v, \mathbf{W}_o$ ) and feed-forward networks ( $\mathbf{W}_{up}, \mathbf{W}_{down}, \mathbf{W}_{gate}$ ). We set the model hyper-parameters based on the optimal configurations from Biderman et al. (2024), employing the decoupled LionW optimizer with a batch size of 192, and training for 8 epochs with a learning rate of  $5e-4$  by default. For RaSA, we set  $k = \max(r/8, 1)$  based on the analysis in § 3.2. More details are provided in appendix C.

### 4.2 MAIN RESULTS

In this section, we compare RaSA and baselines on two challenging domains – code and math.

**Code Generation** Table 1 presents the results on the Humaneval+ test set. We compare RaSA and prior LoRA variants in terms of both efficiency and effectiveness. Although VeRA adds only 1.6M extra parameters for  $r = 1024$ , it results in a training time increase of between 13% and 16% over LoRA with  $r = 32$ . VeRA, however, is the least effective among all the variants due to its extreme strategy for parameter efficiency. Both MoRA and RaSA add a comparable number of additional parameters as LoRA, yet MoRA requires more time due to the use of block diagonal matrices.

In terms of model performance, MoRA shows performance on par with LoRA, aligning with the findings reported in the original paper (Jiang et al., 2024). Our proposed RaSA surpasses all baseline models in nearly all scenarios. Like LoRA, RaSA’s performance improves with rank, and at rank 32, RaSA typically delivers the strongest performance for both the Llama and Mistral models. RaSA achieves maximum Humaneval+ of 59.5% PASS@1 with Llama-3.1-8B.

Table 1: Performance on the code generation task (i.e. Humaneval+). Note that for MoRA and RaSA,  $r$  does not correspond to the effective rank of the update matrix.

r	Method	# Param.	Llama-3.1-8B				Mistral-0.3-7B					
			Time	PASS@1		PASS@10		Time	PASS@1		PASS@10	
				BEST	LAST	BEST	LAST		BEST	LAST	BEST	LAST
1024	VeRA	1.6M	11.3h	48.8	48.8	66.5	64.2	12.5h	42.5	39.5	57.3	54.4
8	LoRA	21.0M	9.6h	56.1	53.0	71.2	68.5	10.7h	42.6	39.7	57.7	54.8
	MoRA	21.0M	12.0h	54.6	52.1	68.4	66.9	13.4h	45.2	38.6	64.4	48.6
	RaSA	21.0M	11.2h	57.9	<b>56.9</b>	<b>72.6</b>	69.6	12.1h	50.0	49.0	66.0	64.2
16	LoRA	41.9M	9.8h	54.5	53.4	68.9	67.6	10.7h	46.0	40.6	61.2	54.9
	MoRA	41.9M	12.7h	56.3	52.9	69.5	65.6	14.0h	43.4	41.0	59.4	56.0
	RaSA	42.0M	11.2h	57.3	56.4	72.1	68.1	12.1h	53.6	51.3	68.5	63.7
32	LoRA	83.9M	10.0h	57.9	<b>56.9</b>	69.8	69.2	10.8h	50.2	44.4	64.4	57.0
	MoRA	83.9M	12.4h	55.6	53.0	69.0	68.3	14.0h	42.2	42.2	56.4	56.0
	RaSA	83.9M	11.5h	<b>59.5</b>	56.2	72.5	<b>71.4</b>	12.5h	<b>55.7</b>	<b>55.7</b>	<b>70.0</b>	<b>65.7</b>

**Mathematical Reasoning** The math results presented in Table 2 are in close alignment with the those of code results. RaSA demonstrates consistent superiority over all baseline models across various configurations. Mistral notably falls short of its Llama counterpart, exhibiting a performance deficit of approximately 8%, which RaSA is capable of narrowing down to 5%. We also observe that directly increasing the hyperparameter  $r$  yields only marginal performance gains, but at the cost of doubling the number of training parameters (see # Param. in Table 1). In contrast, RaSA greatly outperforms LoRA with the same or even fewer parameters ( $RaSA_{r=8} > LoRA_{r=32}$ ). This supports the notion introduced in § 1 that LoRA’s parameters are underutilized. RaSA, on the other hand, improves the utilization of parameters by sharing them across layers.

Table 2: Performance on MATH.

r	Method	Llama-3.1-8B		Mistral-0.3-7B	
		BEST	LAST	BEST	LAST
1024	VeRA	27.4	25.6	19.9	19.4
8	LoRA	28.3	26.7	20.1	19.2
	MoRA	29.2	28.9	21.4	21.4
	RaSA	30.3	29.1	24.3	23.8
16	LoRA	28.8	27.1	20.9	19.5
	MoRA	30.2	26.5	20.5	19.4
	RaSA	31.4	29.8	25.9	25.1
32	LoRA	28.9	27.2	21.8	20.4
	MoRA	28.6	25.8	18.4	18.4
	RaSA	31.7	29.6	26.1	25.1

### 4.3 RASA LEARNS MORE AND FORGETS LESS THAN LORA

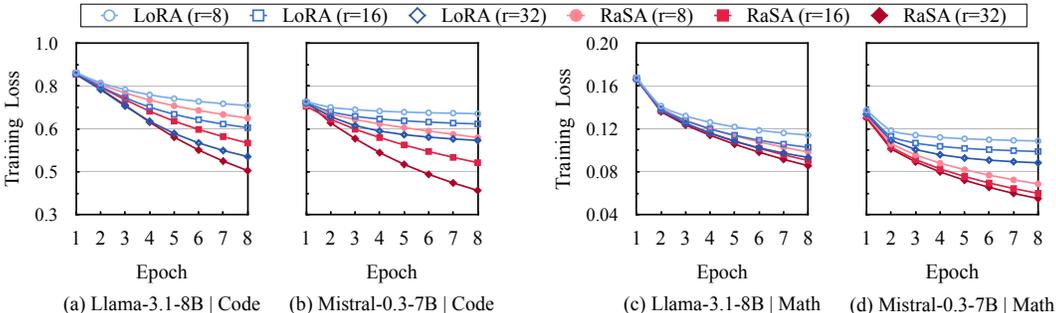


Figure 4: **RaSA learns more and faster than LoRA.** Training curves of LoRA and RaSA with different ranks. RaSA consistently outperforms LoRA with the same rank across models and tasks.

**RaSA learns more and faster than LoRA** Figure 4 illustrates the training curves of the fine-tuning process. Generally, the training losses for both RaSA and LoRA decrease as the rank increases. Notably, RaSA consistently outperforms its LoRA counterpart in terms of both learning effectiveness

and efficiency across all cases, aligning with our empirical analysis presented in Section 3.2. These results collectively underscore the efficacy and universal applicability of the proposed RaSA method. One interesting finding is that RaSA is specifically effective for Mistral: RaSA achieves comparable or potentially superior training outcomes to LoRA with a significantly lower rank requirement of 8, compared to LoRA’s rank of 32.

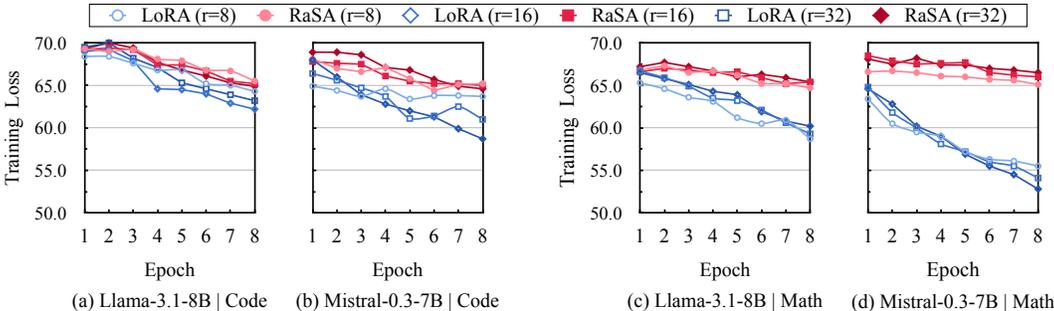


Figure 5: **RaSA forgets less than LoRA.** Y-axis shows the average of prediction accuracy on three benchmarks to evaluate model’s forgetting. Higher prediction accuracy denotes less forgetting.

**RaSA forgets less than LoRA** We follow Biderman et al. (2024) to investigate the extent of forgetting as degradation of base model capabilities. Specifically, we calculate prediction accuracies on the following three benchmarks: (1) HellaSwag (Zellers et al., 2019): inference the most plausible continuation for daily events (70K problems); (2) WinoGrande (Sakaguchi et al., 2019): assesses commonsense reasoning (44K problems); (3) ARC-Challenge (Clark et al., 2018): complex reasoning and understanding of scientific concepts (7.8K problems).

Figure 5 presents the averaged forgetting curves of the three benchmarks, clearly showing that RaSA experiences less forgetting than LoRA, with RaSA’s forgetting levels being less affected by rank changes compared to LoRA. The difference in performance between  $r = 8$  and  $r = 32$  at epoch 8 stands at an average of 2.83% for LoRA and 0.75% for RaSA, indicating a smaller performance variation for RaSA. LoRA is more prone to forgetting in math than code, while RaSA displays greater domain robustness. Specifically, with  $r = 32$ , Mistral scores 58.7% in code and 52.8% in math using LoRA, whereas RaSA shows a reduced performance difference between code (64.6%) and math (66.5%) domains, underscoring RaSA’s robustness.

#### 4.4 SCALING PERFORMANCE ANALYSIS

This section investigates the scaling characteristics of the RaSA approach by varying both the model size and the dataset size to assess its robustness.

**Model Scaling** Initially, we evaluate RaSA’s performance on an expanded scale by examining larger-scale models, including Llama-3.1-70B and Mixtral-8x7B. Due to computational constraints, we employ a rank of  $r = 4$  for both models specifically in the domain of mathematical reasoning. For an equitable comparison, we present results for smaller models configured with  $r = 4$ . Each model is trained over 2 epochs using both LoRA and RaSA techniques, with performance measured in terms of LAST accuracy. The results in Figure 6 reveal that increasing the model size substantially enhances performance for both LoRA and RaSA, across all model types. Noteworthy is the performance of larger Llama and Mistral models using LoRA, achieving MATH accuracies of 40.4% and 32.6%, respectively. These results significantly exceed those of their smaller counterparts under identical configurations and even surpass outcomes from variants with extended training (i.e., 8 epochs). Notably, RaSA consistently outperforms LoRA on these larger-scale models, underscoring RaSA’s robustness in handling models of increased scale.

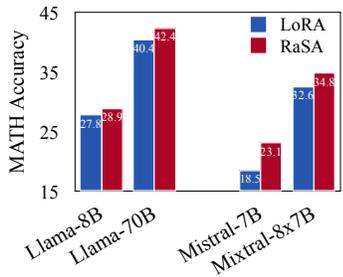


Figure 6: MATH performance of scaled models.

432 **Data Scaling** Subsequently, we explore the influence of training data  
 433 size on RaSA’s performance. We experiment with the Llama-3.1-8B  
 434 model, applying a rank of  $r = 8$  to facilitate efficient training. The  
 435 examination involves random sampling of 25% and 50% instances from  
 436 the SFT data for the mathematics reasoning task. Each model is trained  
 437 over 8 epochs, with performance assessed through the LAST accuracy.  
 438 As illustrated in Figure 7, LoRA’s performance seems contingent on the  
 439 volume of training data, with no noticeable improvement when data is  
 440 increased from 25% to 50%. This finding is consistent with the results  
 441 in Biderman et al. (2024). In contrast, RaSA demonstrates a remarkable  
 442 ability to enhance performance with an increase in training data volume.  
 443 Impressively, with just 25% of the training data, RaSA outperforms  
 444 LaSA even when the latter utilizes the entire dataset, highlighting  
 445 RaSA’s exceptional efficiency in leveraging training data for performance improvement.

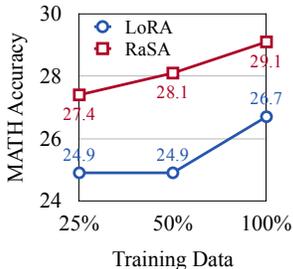


Figure 7: MATH performance of scaled data.

446  
 447 **5 RELATED WORK**

448 **Parameter-Efficient Fine-Tuning (PEFT)** PEFT methods aim to minimize the number of trainable  
 449 parameters needed for fine-tuning large models, thus reducing memory and computational require-  
 450 ments. Pioneering methods include adapter-based (Houlsby et al., 2019) and prompt-based (Lester  
 451 et al., 2021; Li & Liang, 2021) that introduce additional tunable adapter or prefix tokens to enable  
 452 efficient fine-tuning while keeping the original model parameters fixed. However, these approaches  
 453 can slow down inference speed due to the extra components introduced. LoRA overcomes this  
 454 drawback by introducing low-rank matrices directly into the weight update process during fine-tuning,  
 455 effectively reducing trainable parameters without increasing inference latency. Due to its robust  
 456 performance, LoRA and its variants have been widely used to adapt LLMs for specific tasks (Yu  
 457 et al., 2024; Xu et al., 2023; Biderman et al., 2024; Chen et al., 2024; Meng et al., 2024; yang Liu  
 458 et al., 2024). [Benedek & Wolf \(2024\)](#) and [Zhang et al. \(2023\)](#) show that the number of ranks required  
 459 for each parameter matrix across the model’s layers is not uniform. Therefore, they propose dynam-  
 460 ically assigning ranks based on the importance of parameters during training. These rank-allocating  
 461 approaches typically involve real-time estimation of parameter importance and pruning during the  
 462 training process. In contrast, RaSA uses a shared rank pool combined with layer-specific weighting,  
 463 eliminating the need for complex importance estimation or pruning. Biderman et al. (2024) conduct  
 464 a comprehensive empirical study on LoRA, and reveal that while LoRA still lags behind FFT, it  
 465 exhibits less catastrophic forgetting. We show that our proposed RaSA forgets even less than LoRA,  
 466 and learns more and faster.

467 **Parameter Redundancy of LoRA** Although LoRA has significantly reduced the number of  
 468 trainable parameters, recent research suggest that it is possible to further minimize these parameters  
 469 without compromising performance. Kopiczko et al. (2024) achieve a 99% reduction in LoRA  
 470 parameters by fully sharing a pair of low-rank, frozen random matrices across all layers, adjusted with  
 471 learnable scaling vectors. Koohpayegani et al. (2024) propose learning linear combinations of a set  
 472 of random matrix bases, while Li et al. (2024) push this further by replacing the matrix bases with a  
 473 vector bank. Song et al. (2024) and Renduchintala et al. (2024) explore the effects of different sharing  
 474 and selective fine-tuning strategies. By sharing parameter spaces, Brüel-Gabrielsson et al. (2024)  
 475 compress 1,000 LoRAs trained from different task, enabling more efficient serving. These findings  
 476 collectively suggest that LoRA’s parameter has not been fully utilized and that different LoRAs  
 477 exhibit similarities across layers, modules, and even different tasks. Rather than focusing on extreme  
 478 parameter reduction, this work aims to maintain the same parameter count while exploring how inter-  
 479 layer sharing can enhance parameter utilization. We theoretically and empirically demonstrate that  
 480 sharing ranks across layers leads to lower reconstruction error and thus better expressive capacity.

481  
 482 **6 CONCLUSION**

483  
 484 In this study, we introduced RaSA, a novel extension to LoRA through an innovative partial rank  
 485 sharing across layers. RaSA maintains the parameter efficiency and seamless integration into existing  
 models characteristic of LoRA while substantially increasing the model’s expressiveness. Through

486 theoretical analysis, we established RaSA’s superior capability in matrix reconstruction compared  
487 to traditional LoRA, underpinning its improved performance in downstream tasks. Empirical re-  
488 sults on complex tasks such as code generation and mathematical reasoning have demonstrated its  
489 effectiveness over LoRA in high-demand scenarios. Future research directions may explore further  
490 optimization of rank-sharing schemes and the potential of RaSA in a broader range of applications,  
491 paving the way for the development of even more powerful and efficient PEFT strategies.

#### 492 REPRODUCIBILITY STATEMENT

493 We are committed to ensuring the reproducibility of our results. All code, data, and scripts used for  
494 our experiments are available at <https://anonymous.4open.science/r/RaSA-ICLR-0E25>. All  
495 datasets and models used in our work are publicly accessible. Detailed training and data details are  
496 provided in appendix C.  
497

#### 498 REFERENCES

- 501 Loubna Ben Allal, Niklas Muennighoff, Logesh Kumar Umapathi, Ben Lipkin, and Leandro von  
502 Werra. A framework for the evaluation of code generation models. [https://github.com/](https://github.com/bigcode-project/bigcode-evaluation-harness)  
503 [bigcode-project/bigcode-evaluation-harness](https://github.com/bigcode-project/bigcode-evaluation-harness), 2022.
- 504 Nadav Benedek and Lior Wolf. PRILoRA: Pruned and rank-increasing low-rank adaptation. In  
505 Yvette Graham and Matthew Purver (eds.), *Findings of the Association for Computational Linguis-*  
506 *tics: EACL 2024*, pp. 252–263, St. Julian’s, Malta, March 2024. Association for Computational  
507 Linguistics. URL <https://aclanthology.org/2024.findings-eacl.18>.
- 508 Dan Biderman, Jose Gonzalez Ortiz, Jacob Portes, Mansheej Paul, Philip Greengard, Connor Jennings,  
509 Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. Lora learns less and forgets less.  
510 *arXiv preprint arXiv:2405.09673*, 2024.
- 511 Rickard Brüel-Gabrielsson, Jiacheng Zhu, Onkar Bhardwaj, Leshem Choshen, Kristjan Greenewald,  
512 Mikhail Yurochkin, and Justin Solomon. Compress then serve: Serving thousands of lora adapters  
513 with little overhead, 2024. URL <https://arxiv.org/abs/2407.00066>.
- 514 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared  
515 Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large  
516 language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- 517 Yukang Chen, Shengju Qian, Haotian Tang, Xin Lai, Zhijian Liu, Song Han, and Jiaya Jia. Long-  
518 gloRA: Efficient fine-tuning of long-context large language models. In *The Twelfth International*  
519 *Conference on Learning Representations, 2024*. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=6PmJoRfdaK)  
520 [6PmJoRfdaK](https://openreview.net/forum?id=6PmJoRfdaK).
- 521 Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and  
522 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.  
523 *arXiv:1803.05457v1*, 2018.
- 524 Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser,  
525 Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. Training verifiers to solve  
526 math word problems. *arXiv preprint arXiv:2110.14168*, 2021.
- 527 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha  
528 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.  
529 *arXiv preprint arXiv:2407.21783*, 2024.
- 530 Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychome-*  
531 *trika*, 1(3):211–218, 1936.
- 532 Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster,  
533 Laurence Golding, Jeffrey Hsu, Alain Le Noac’h, Haonan Li, Kyle McDonell, Niklas Muennighoff,  
534 Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika,  
535 Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot  
536 language model evaluation, 07 2024. URL <https://zenodo.org/records/12608602>.

- 540 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing  
541 human-level performance on imagenet classification. In *2015 IEEE International Conference on*  
542 *Computer Vision (ICCV)*, pp. 1026–1034, 2015. doi: 10.1109/ICCV.2015.123.
- 543  
544 Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song,  
545 and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In  
546 *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*  
547 *(Round 2)*, 2021. URL <https://openreview.net/forum?id=7Bywt2mQsCe>.
- 548 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe,  
549 Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning  
550 for NLP. In Kamalika Chaudhuri and Ruslan Salakhutdinov (eds.), *Proceedings of the 36th*  
551 *International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*  
552 *Research*, pp. 2790–2799. PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/houlsby19a.html>.
- 553  
554 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
555 Weizhu Chen. LoRA: Low-rank adaptation of large language models. In *International Conference*  
556 *on Learning Representations*, 2022. URL <https://openreview.net/forum?id=nZeVKeeFYf9>.
- 557  
558 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
559 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.  
560 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 561 Ting Jiang, Shaohan Huang, Shengyue Luo, Zihan Zhang, Haizhen Huang, Furu Wei, Weiwei Deng,  
562 Feng Sun, Qi Zhang, Deqing Wang, et al. Mora: High-rank updating for parameter-efficient  
563 fine-tuning. *arXiv preprint arXiv:2405.12130*, 2024.
- 564  
565 Soroush Abbasi Koohpayegani, Navaneet K L, Parsa Nooralinejad, Soheil Kolouri, and Hamed  
566 Pirsiavash. NOLA: Compressing lora using linear combination of random basis. In *The Twelfth*  
567 *International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=TjfXcDgvzk>.
- 568  
569 Dawid Jan Kopiczko, Tijmen Blankevoort, and Yuki M Asano. VeRA: Vector-based random matrix  
570 adaptation. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
571 <https://openreview.net/forum?id=NjNfLdxr3A>.
- 572  
573 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient  
574 prompt tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-  
575 tau Yih (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Lan-*  
576 *guage Processing*, pp. 3045–3059, Online and Punta Cana, Dominican Republic, November  
577 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.243. URL  
<https://aclanthology.org/2021.emnlp-main.243>.
- 578  
579 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation.  
580 In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.), *Proceedings of the 59th*  
581 *Annual Meeting of the Association for Computational Linguistics and the 11th International Joint*  
582 *Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 4582–4597, Online,  
583 August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.353.  
URL <https://aclanthology.org/2021.acl-long.353>.
- 584  
585 Yang Li, Shaobo Han, and Shihao Ji. Vb-lora: Extreme parameter efficient fine-tuning with vector  
586 banks. *arXiv preprint arXiv:2405.15179*, 2024.
- 587  
588 Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated  
589 by chatGPT really correct? rigorous evaluation of large language models for code generation.  
590 In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=1qvx610Cu7>.
- 591  
592 Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing  
593 Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with  
evol-instruct. In *The Twelfth International Conference on Learning Representations*, 2024. URL  
<https://openreview.net/forum?id=UnUwSIgK5W>.

- 594 Fanxu Meng, Zhaohui Wang, and Muhan Zhang. Pissa: Principal singular values and singular vectors  
595 adaptation of large language models. *arXiv preprint arXiv:2404.02948*, 2024.  
596
- 597 Adithya Renduchintala, Tugrul Konuk, and Oleksii Kuchaiev. Tied-LoRA: Enhancing parameter  
598 efficiency of LoRA with weight tying. In Kevin Duh, Helena Gomez, and Steven Bethard  
599 (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association  
600 for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp.  
601 8694–8705, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi:  
602 10.18653/v1/2024.naacl-long.481. URL <https://aclanthology.org/2024.naacl-long.481>.
- 603 Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An  
604 adversarial winograd schema challenge at scale. *arXiv preprint arXiv:1907.10641*, 2019.  
605
- 606 Yurun Song, Junchen Zhao, Ian G Harris, and Sangeetha Abdu Jyothi. Sharelora: Parameter efficient  
607 and robust large language model fine-tuning via shared low-rank adaptation. *arXiv preprint  
608 arXiv:2406.10785*, 2024.
- 609 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N  
610 Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon,  
611 U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett  
612 (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Asso-  
613 ciates, Inc., 2017. URL [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/  
614 3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf).
- 615 Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering  
616 code generation with OSS-instruct. In Ruslan Salakhutdinov, Zico Kolter, Katherine Heller, Adrian  
617 Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceedings of the 41st  
618 International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning  
619 Research*, pp. 52632–52657. PMLR, 21–27 Jul 2024. URL [https://proceedings.mlr.press/  
620 v235/wei24h.html](https://proceedings.mlr.press/v235/wei24h.html).
- 621 Haoran Xu, Young Jin Kim, Amr Sharaf, and Hany Hassan Awadalla. A paradigm shift in machine  
622 translation: Boosting translation performance of large language models, 2023.  
623
- 624 Shih yang Liu, Chien-Yi Wang, Hongxu Yin, Pavlo Molchanov, Yu-Chiang Frank Wang, Kwang-Ting  
625 Cheng, and Min-Hung Chen. DoRA: Weight-decomposed low-rank adaptation. In *Forty-first  
626 International Conference on Machine Learning*, 2024. URL [https://openreview.net/forum?  
627 id=3d5CIRG1n2](https://openreview.net/forum?id=3d5CIRG1n2).
- 628 Longhui Yu, Weisen Jiang, Han Shi, Jincheng YU, Zhengying Liu, Yu Zhang, James Kwok, Zhenguo  
629 Li, Adrian Weller, and Weiyang Liu. Metamath: Bootstrap your own mathematical questions for  
630 large language models. In *The Twelfth International Conference on Learning Representations*,  
631 2024. URL <https://openreview.net/forum?id=N8N0hgNDRt>.
- 632 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. HellaSwag: Can a  
633 machine really finish your sentence? In Anna Korhonen, David Traum, and Lluís Màrquez  
634 (eds.), *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*,  
635 pp. 4791–4800, Florence, Italy, July 2019. Association for Computational Linguistics. doi:  
636 10.18653/v1/P19-1472. URL <https://aclanthology.org/P19-1472>.
- 637
- 638 Yuchen Zeng and Kangwook Lee. The expressive power of low-rank adaptation. In *The Twelfth  
639 International Conference on Learning Representations*, 2024. URL [https://openreview.net/  
640 forum?id=1ikXVjmh3E](https://openreview.net/forum?id=1ikXVjmh3E).
- 641 Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen,  
642 and Tuo Zhao. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh  
643 International Conference on Learning Representations*, 2023. URL [https://openreview.net/  
644 forum?id=lq62uWRJjiY](https://openreview.net/forum?id=lq62uWRJjiY).
- 645  
646  
647

## A COORDINATE DESCENT EXPERIMENT

This section details the derivation of the coordinate descent experiment discussed in § 3.2, inspired by Brüel-Gabrielsson et al. (2024).

Given the parameters of RaSA:  $\{\tilde{\mathbf{B}}_i, \tilde{\mathbf{A}}_i, \mathbf{B}_S, \tilde{\mathbf{D}}_i, \mathbf{A}_S\}_{i \in [L]}$ , the reconstruction error of RaSA is defined as:

$$E = \sum_{i=1}^L \|\mathbf{M}_i - (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S)\|_F^2. \quad (17)$$

Clearly,  $\tilde{\mathbf{B}}_i$  and  $\tilde{\mathbf{A}}_i$  are independent across layers. By applying the Eckart–Young–Mirsky theorem (Eckart & Young, 1936), we first compute the SVD of the residual matrix:

$$\text{SVD}(\mathbf{M}_i - \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S) = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T. \quad (18)$$

Therefore, the update rules for  $\tilde{\mathbf{B}}_i$  and  $\tilde{\mathbf{A}}_i$  are:

$$\begin{aligned} \tilde{\mathbf{B}}_i &= \mathbf{U}_{[:, :r-k]} \mathbf{\Sigma}_{[r-k, :r-k]}^{\frac{1}{2}}, \\ \tilde{\mathbf{A}}_i &= \mathbf{\Sigma}_{[r-k, :r-k]}^{\frac{1}{2}} \mathbf{V}_{[:, :r-k]}^T. \end{aligned} \quad (19)$$

Let the low-rank decomposition of  $\mathbf{M}_i$  be:  $\mathbf{M}_i = \hat{\mathbf{B}}_i \hat{\mathbf{A}}_i$ , where  $\hat{\mathbf{B}}_i \in \mathbb{R}^{b \times R}$  and  $\hat{\mathbf{A}}_i \in \mathbb{R}^{R \times a}$ . Next, we compute the following gradients:

$$\nabla_{\mathbf{B}_S} E = \sum_{i=1}^L -2 \left( \hat{\mathbf{B}}_i \hat{\mathbf{A}}_i - (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S) \right) \mathbf{A}_S^T \mathbf{D}_i^T, \quad (20)$$

$$\nabla_{\mathbf{A}_S} E = \sum_{i=1}^L -2 \left( \hat{\mathbf{B}}_i \hat{\mathbf{A}}_i - (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S) \right) \mathbf{B}_S \mathbf{D}_i, \quad (21)$$

$$\nabla_{\mathbf{D}_i} E = -2 \mathbf{B}_S^T \left( \hat{\mathbf{B}}_i \hat{\mathbf{A}}_i - (\tilde{\mathbf{B}}_i \tilde{\mathbf{A}}_i + \mathbf{B}_S \mathbf{D}_i \mathbf{A}_S) \right) \mathbf{A}_S^T, \quad (22)$$

$$\nabla_{\text{diag}(\mathbf{D}_i)} E = \text{diag}(\nabla_{\mathbf{D}_i} E). \quad (23)$$

By setting these gradients to zeros, we obtain the following update rules:

$$\mathbf{B}_S = \left( \sum_{i=1}^L \begin{bmatrix} \hat{\mathbf{B}}_i & -\tilde{\mathbf{B}}_i \\ \hat{\mathbf{A}}_i & \tilde{\mathbf{A}}_i \end{bmatrix} \mathbf{A}_S^T \mathbf{D}_i^T \right) \left( \sum_{i=1}^L \mathbf{D}_i \mathbf{A}_S \mathbf{A}_S^T \mathbf{D}_i^T \right)^{-1}, \quad (24)$$

$$\mathbf{A}_S = \left[ \left( \sum_{i=1}^L \left( \begin{bmatrix} \hat{\mathbf{B}}_i & -\tilde{\mathbf{B}}_i \\ \hat{\mathbf{A}}_i & \tilde{\mathbf{A}}_i \end{bmatrix} \right)^T \mathbf{B}_S \mathbf{D}_i \right) \left( \sum_{i=1}^L \mathbf{D}_i^T \mathbf{B}_S^T \mathbf{B}_S \mathbf{D}_i \right)^{-1} \right]^T, \quad (25)$$

$$\text{diag}(\mathbf{D}_i) = (\mathbf{B}_S^T \mathbf{B}_S \circ \mathbf{A}_S \mathbf{A}_S^T)^{-1} \left( \mathbf{B}_S^T \begin{bmatrix} \hat{\mathbf{B}}_i & -\tilde{\mathbf{B}}_i \end{bmatrix} \circ \mathbf{A}_S \begin{bmatrix} \hat{\mathbf{A}}_i \\ \tilde{\mathbf{A}}_i \end{bmatrix}^T \right) \mathbf{1}. \quad (26)$$

In coordinate descent, we iteratively apply Equations (19) and (24) to (26) until convergence.

## 702 B PROMPT TEMPLATES

703

704

705

706

Below is an instruction that describes a task. Write a response that appropriately completes the request.

707

708

709

### Instruction:

{QUESTION}

710

711

712

### Response:

Let's think step by step.

713

714

Figure 8: Evaluation prompt for mathematical reasoning.

715

716

717

Below is an instruction that describes a task. Write a response that appropriately completes the request.

718

719

720

### Instruction:

{QUESTION}

721

722

723

### Response:

{IMPORT SECTION}

724

725

{FUNCTION SIGNATURE}

{DOCSTRING}

726

727

728

729

Figure 9: Evaluation prompt for code generation.

730

731

## C TRAINING AND DATA DETAILS

732

733

734

735

736

737

738

739

740

741

742

**Training** We mostly aligned our training configurations with the optimal configurations from Biderman et al. (2024). For LoRA, we used the decoupled LionW optimizer with a batch size of 192, training for 8 epochs with a learning rate of  $5e-4$ . A cosine learning rate scheduler was applied, with the first 10% of training steps used for warmup, and weight decay set to zero. Training and evaluation were conducted using bfloat16 precision. While Biderman et al. (2024) set  $\alpha = 32$  for both math and code tasks, we reduced  $\alpha$  to 8 for the math task due to convergence issues observed with the Mistral model when using  $\alpha = 32$ . RaSA training fully inherits all hyper-parameters from LoRA training. For MoRA, we used a learning rate of  $3e-4$ , as reported in the original work (Jiang et al., 2024). For VeRA, following the original paper, we set the learning rate to 10 times that of LoRA, resulting in  $5e-3$  (Kopiczko et al., 2024). All experiments for the 7-8B models were conducted on 1 node  $\times$  8  $\times$  A100-40G GPUs. For the 70B and MoE models, we used 8 nodes.

743

744

745

746

**Data** During training, we grouped data by length, which significantly accelerated the training process. All math training data ends with “The answer is: {ANSWER}”, helping answer extraction during evaluation.

747

748

## D ADDING MORE BASELINES

749

750

751

752

753

754

755

Table 3 is adding more baseline in Table 2. We only present partial results since time constraint and will complete full results in the final version.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

Table 3: Performance on MATH (adding more baselines in Table 2).

r	Method	# Trainable Param.	# Extra Param.	Llama-3.1-8B		Mistral-0.3-7B	
				BEST	LAST	BEST	LAST
–	FFT	7-8B	00.0M	xx.x	xx.x	28.1	26.6
1024	VeRA	1.6M	00.0M	27.4	25.6	19.9	19.4
8	LoRA	21.0M	00.0M	28.3	26.7	20.1	19.2
	MoRA	21.0M	00.0M	29.2	28.9	21.4	21.4
	OLoRA	21.0M	00.0M	xx.x	xx.x	22.5	22.5
	AdaLoRA	31.5M	63.0M	xx.x	xx.x	22.5	21.6
	PriLoRA	21.3M	10.7M	xx.x	xx.x	22.3	22.3
	RaSA	21.0M	00.0M	30.3	29.1	24.3	23.8
16	LoRA	41.9M	00.0M	28.8	27.1	20.9	19.5
	MoRA	41.9M	00.0M	30.2	26.5	20.5	19.4
	OLoRA	41.9M	00.0M	xx.x	xx.x	22.5	22.2
	AdaLoRA	62.9M	125.8M	xx.x	xx.x	23.5	23.2
	PriLoRA	42.6M	21.3M	xx.x	xx.x	22.7	21.6
	RaSA	42.0M	00.0M	31.4	29.8	25.9	25.1
32	LoRA	83.9M	00.0M	28.9	27.2	21.8	20.4
	MoRA	83.9M	00.0M	28.6	25.8	18.4	18.4
	OLoRA	83.9M	00.0M	xx.x	xx.x	xx.x	xx.x
	AdaLoRA	125.9M	251.8M	xx.x	xx.x	xx.x	xx.x
	PriLoRA	85.2M	42.6M	xx.x	xx.x	xx.x	xx.x
	RaSA	83.9M	00.0M	31.7	29.6	26.1	25.1